```python
from datetime import datetime, timedelta
from dateutil.parser import parse
import random
import pandas as pd
import numpy as np

import sklearn
from sklearn.metrics import silhouette_score
from sklearn import preprocessing

import matplotlib.pyplot as plt
from matplotlib import cm

from pyclustering.cluster.kmeans import kmeans
from pyclustering.utils.metric import type_metric, distance_metric
from pyclustering.utils.metric \
    import euclidean_distance, manhattan_distance, chebyshev_distance, \
    minkowski_distance
from pyclustering.cluster.center_initializer import kmeans_plusplus_initializer


class PycClusteringPeople:
    df_corona = None  # initial loaded data
    added_column_list = []  # Columns added by calculation within class

    sse_list = []  # list for SSE of each cluster
    sil_score_list = []  # list for Silhouette score of clustering result
    centroids_coord_list = []  # list for storing coordinates of centroids

    # member variable for custom distance function
    weight_list = []  # weight values list

    # TODO: Discard below 4 values
    num_of_data = 0
    data_cal_count = 0
    num_of_cent = 0
    cent_cal_count = 0

    def __init__(self, file_path):
        self.load_data(file_path)
        self.compute_severity(parse('2020 7 3'))

    def load_data(self, file_path):
        """
         method to load .csv file
         :param file_path: string, the path of file
         :return:
        """
        self.df_corona = pd.read_csv(file_path)

    def compute_severity(self, base_date):
        """
         method to preprocess the data for distance function
         :return: None
        """
        col_num = len(self.df_corona)  # the number of rows from loaded data
        today = datetime.now().date()  # date of today, YEAR-MONTH-DAY

        # selecting specific column to compute 'severity'

        severity_list = []  # list for storing severity result

        for i in range(col_num):
            incur_date_col = self.df_corona['Incurred Date']
            status = self.df_corona['Covid Status']
            severity = 0  # default is healthy. 0.
            if status[i] == 'Contacted':  # contacted person?
```

```python
                        # formula for contacted person:
                        #   x = 1 - ((today's date) - (infected date)) * 0.05)
                        elapsed_days = (base_date.date() - parse(incur_date_col[i]).date()).days
                        severity = (1 - (elapsed_days * 0.05)) * 0.5

                    elif status[i] == 'Confirmed':  # confirmed person?
                        # formula for confirmed person:
                        #   x = (1 - ((today's date) - (infected date)) * 0.05)) / 2
                        elapsed_days = (base_date.date() - parse(incur_date_col[i]).date()).days
                        severity = 1 - (elapsed_days * 0.05)

                    # add the value to the list
                    # and rounding to solve floating-point problems
                    severity_list.append(round(severity, 4))
                self.df_corona["Severity"] = severity_list
                self.added_column_list.append("Severity")

    def initialize_random_centroid(self, num_centroid, is_0_1_normalized=True):
        """
        A function that generates a centroid of random coordinates-
        -as many as the number of clusters received.
        :param num_centroid: int, the number of centroids
        :param is_0_1_normalized: boolean, Whether the feature is normalized
        :return: the random coordinates of centroids list
        """
        if is_0_1_normalized:  # are all columns normalized to 0-1?
            # for i in num_centroid:
            return [[random.uniform(0, 1), random.uniform(0, 1)] for _ in range(
                num_centroid)]

        else:  # there is an unnormalized column
            pass

    def pyc_cluster_kmeans(self,
                           target_col_name_list,
                           num_cluster,
                           weight_list,
                           distance_function):
        """
        function to cluster data
        :param target_col_name_list: list, name list to extract column as feature
        :param num_cluster: int, the number of clusters
        :param weight_list: list, weight list of features
        :param distance_function: string, the abbreviation of distance function
        :return: list, clustered result
        """
        self.weight_list = weight_list

        # my_distance_function = lambda p1, p2: p1[0] + p2[0] + 2
        if distance_function == 'eu':
            # metric = distance_metric(type_metric.EUCLIDEAN)
            metric = euclidean_distance
        elif distance_function == 'ma':
            # metric = distance_metric(type_metric.MANHATTAN)
            metric = manhattan_distance
        elif distance_function == 'mi':
            # metric = distance_metric(type_metric.MINKOWSKI)
            metric = minkowski_distance
        elif distance_function == 'c_eu':
            metric = distance_metric(type_metric.USER_DEFINED, func=self.
                weighted_euclidean_distance)

        target_data = self.df_corona.loc[:, target_col_name_list]  # To select
        required data
        if "Age" in target_col_name_list:  # scaling only "Age" column
            age = target_data.loc[:, "Age"].values.reshape(-1, 1)
            scaler = preprocessing.MinMaxScaler()
```

```python
131                        # data = scaler.fit_transform(data)    # To scale data from 0 to 1
132                        target_data.loc[:, "Age"] = scaler.fit_transform(age)
133
134                # set the number of data and centroids
135                self.data_cal_count = num_cluster*len(target_data)
136                self.num_of_data = num_cluster*len(target_data)
137                self.cent_cal_count = 1
138                self.num_of_cent = 1
139
140                # initializing centroids
141                # initial_centers = kmeans_plusplus_initializer(target_data,
                   num_cluster).initialize()
142                # initial_centers = self.initialize_random_centroid(num_cluster)
143                initial_centers = [[i*0.1, i*0.1] for i in range(num_cluster)]
144
145                kmeans_instance = kmeans(target_data, initial_centers, metric=metric)
146
147                kmeans_instance.process()
148                clustered_list = kmeans_instance.get_clusters()
149                clustered_list = self.pyc_result_to_column(clustered_list, len(target_data))
150
151                # add the column
152                self.df_corona['Cluster ID'] = clustered_list
153
154                # storing the coordinates of centroids
155                self.centroids_coord_list.append(kmeans_instance.get_centers())
156
157                # storing SSE(Sum of Squared Errors)
158                self.sse_list.append(kmeans_instance.get_total_wce())
159
160                # strong Silhouette Score
161                self.sil_score_list.append(silhouette_score(target_data, clustered_list))
162
163                return clustered_list
164
165        def pyc_result_to_column(self, pyc_cluster_result, people_num):
166            """
167             function to change shape of Pyclustering to pandas
168             :param pyc_cluster_result: nd list, result of clustering using Pyclusterin
169             :param people_num: int, the number of people(data)
170             :return: list, re-shaped list
171            """
172            clustered_list = [0 for _ in range(people_num)]
173
174            cluster_id = 0
175            for id_list_of_a_cluster in pyc_cluster_result:
176                for idx in id_list_of_a_cluster:
177                    clustered_list[idx] = cluster_id
178                cluster_id += 1
179
180            return clustered_list
181
182        def weighted_euclidean_distance(self, point1, point2):
183            """
184             custom distance function
185             :param point1: list, list of feature values or coordinates list of centroid
186             :param point2: coordinates list of centroid
187             :return: distance between point1 and point2
188            """
189            distance = 0.0  # distance between point1 and point2
190
191            # TODO: condition is should be changed to type comparison
192            if self.data_cal_count > 0:  # when calculating the distance of two coordinates
193                # point 1 is data.
194                # point 2 is centroid
195                self.data_cal_count -= 1
196                for weight, p1_coord, p2_coord in zip(self.weight_list, point1, point2):
```

```python
197                    distance += weight * (p1_coord - p2_coord) ** 2.0
198
199            else:  # when updating Centroid
200                # point 1 and 2 are centroid
201                self.cent_cal_count -= 1
202                if self.cent_cal_count == 0:
203                    self.data_cal_count = self.num_of_data
204                    self.cent_cal_count = self.num_of_cent
205
206                for prev_cent, curr_cent in zip(point1, point2):
207                    for weight, pc_coord, cc_coord in zip(self.weight_list, prev_cent,
                        curr_cent):
208                        distance += weight * (pc_coord - cc_coord) ** 2.0
209
210            return distance ** 0.5
211
212        def display_clustering_result(self,
213                                      num_cluster,
214                                      cluster_idx_list,
215                                      cluster_predicted_list):
216            """
217            function to display clustering result on console as tabular type
218            :param num_cluster: int, the number of cluster
219            :param cluster_idx_list: list, cluster index list, ie. [2, 3, 4, 5, 6]
220            :param cluster_predicted_list:
221            :return:
222            """
223
224            severity_list = self.df_corona["Severity"].values.tolist()
225            age_list = self.df_corona["Age"].values.tolist()
226            if type(cluster_predicted_list) != list:
227                cluster_predicted_list = cluster_predicted_list.tolist()
228            people_num_of_a_cluster_list = []
229            avg_age_of_a_cluster_list = []
230            avg_severity_of_a_cluster_list = []
231
232            print(f"Number of Clusters: {len(cluster_idx_list)}")
233
234            for cluster_idx in cluster_idx_list:  # 1 cluster
235                num_people = cluster_predicted_list.count(cluster_idx)
236                id_target_data_tuple_list = []
237                target_severity_list = []
238                target_age_list = []
239
240                for person_idx in range(len(cluster_predicted_list)):
241                    if cluster_idx == cluster_predicted_list[person_idx]:
242                        target_severity_list.append(severity_list[person_idx])
243                        target_age_list.append(age_list[person_idx])
244                        id_target_data_tuple_list.append((
245                            person_idx+1,  # [0] of tuple is id
246                            age_list[person_idx],  # [1] of tuple is age
247                            round(severity_list[person_idx], 2)))  # [2] of tuple is
                            severity
248
249                people_num_of_a_cluster_list.append(num_people)
250
251                print(f"\tCluster {cluster_idx}:")
252                print(f"\t\tNumber of People: {num_people}")
253                # print(f"\t\t\t{'ID':<4}{'Age':<4}{'Severity Value'}")
254                # for person_in_cluster in id_target_data_tuple_list:
255                #     print(f"\t\t\t{person_in_cluster[0]:<4}"
256                #           f"{person_in_cluster[1]:<4}"
257                #           f"{person_in_cluster[2]}")
258                print(f"\t\tMinimum of Age values: {min(target_age_list)}")
259                print(f"\t\tMaximum of Age values: {max(target_age_list)}")
260                print(f"\t\tAverage of Age values: "
261                      f"{round(sum(target_age_list) / len(id_target_data_tuple_list), 2)}")
```

```python
262                print(f"\t\tMinimum of Severity values: {min(target_severity_list)}")
263                print(f"\t\tMaximum of Severity values: {max(target_severity_list)}")
264                print(f"\t\tAverage of Severity values: "
265                      f"{round(sum(target_severity_list) / len(id_target_data_tuple_list),
266                      2)}")
                   print(f"\t\tThe Coordinates of Centroid:")
267                coords = self.centroids_coord_list[num_cluster-2][cluster_idx]
268                print(f"\t\t\tX1 (Severity): {round(coords[0], 2)}")
269                print(f"\t\t\tX2 (Age): {round(coords[1], 2)}")
270                try:
271                    avg_age_of_a_cluster_list.append(
272                        round(sum(target_age_list) / len(id_target_data_tuple_list), 2))
273                except ZeroDivisionError:
274                    avg_age_of_a_cluster_list.append(0)
275
276                try:
277                    avg_severity_of_a_cluster_list.append(
278                        round(sum(target_severity_list) / len(id_target_data_tuple_list), 2
                            ))
279                except ZeroDivisionError:
280                    avg_severity_of_a_cluster_list.append(0)
281
282                print()  # float 1 line
283            self.display_summary_table(people_num_of_a_cluster_list,
284                                       avg_age_of_a_cluster_list,
285                                       avg_severity_of_a_cluster_list)
286        print()  # float 1 line
287
288    def data_to_csv(self, num_cluster):
289        """
290         function to save data as .csv file
291         :param num_cluster: int, the number of cluster.
292         :return: None
293        """
294        temp_df = self.df_corona.__deepcopy__()
295
296        file_name = f"clustered_corona_data_k={num_cluster}_" \
297                    f"{'Severity_Age'}_{''.join(str(self.weight_list))}.csv"
298        temp_df.to_csv(file_name, encoding='utf-8-sig')
299
300    def display_load_data(self):
301        """
302         function to display data
303         :return: None
304        """
305        print(f"Total number of People: {len(self.df_corona)}")
306        print(f"{'ID':<4}"
307              f"{'Age':<4}"
308              f"{'Covid Status':<13}"
309              f"{'Severity':<9}"
310              f"{'Address':<10}")
311        for i in range(len(self.df_corona)):
312            print(f"{self.df_corona['ID'][i]:<4}"
313                  f"{self.df_corona['Age'][i]:<4}"
314                  f"{self.df_corona['Covid Status'][i]:<13}"
315                  f"{round(self.df_corona['Severity'][i], 3):<9}"
316                  f"{self.df_corona['Address'][i].split()[0]:<10}"
317                  )
318        print()  # float 1 line
319        grouped_status = self.df_corona['Severity'].groupby(self.df_corona['Covid
                Status'])
320
321        print(f"Number of healthy people: {grouped_status.count()['Healthy']}")
322        print(f"Number of contacted people: {grouped_status.count()['Contacted']}")
323        print(f"Number of confirmed people: {grouped_status.count()['Confirmed']}")
324
325        print(f"Average Severity of contacted people: "
```

```python
326                         f"{round(grouped_status.mean()['Contacted'], 2)}")
327             print(f"Average Severity of confirmed people: "
328                         f"{round(grouped_status.mean()['Confirmed'], 2)}")
329             print()  # float 1 line
330
331     def plot_data(self, num_cluster):
332         """
333          To plot result
334          :return:
335          """
336         groups = self.df_corona.groupby("Cluster ID")
337         fig, ax = plt.subplots()
338         for name, group in groups:
339             ax.plot(group.Severity, group.Age, marker='o', linestyle="", label=name)
340         ax.legend(fontsize=12)
341         plt.title("Result of Clustering (K="+str(num_cluster)+", Weight="+str(self.
                weight_list)+")")
342         plt.xlabel("Severity")
343         plt.ylabel("Age")
344         # plt.show()
345         fig.savefig("./Cluster_Result_Plotting/cluster_result_"+str(num_cluster)+"_"+
                str(self.weight_list)+".png", dpi=300)
346         plt.close()
347
348     def display_summary_table(self,
349                                 people_num_of_a_cluster_list,
350                                 avg_age_of_cluster_list,
351                                 avg_severity_of_cluster_list):
352         """
353          function to display the data as tabular summary
354          :param people_num_of_a_cluster_list: list, the number of people in a cluster
355          :param avg_age_of_cluster_list: list,
356          :param avg_severity_of_cluster_list:
357          :return:
358          """
359         len_id = 17
360         len_p_num = 11
361         len_age = 13
362         len_sev = 15
363         len_sum = len_id+len_p_num+len_age+len_sev
364
365         # top row
366         print(f"\t{'-'*(len_sum+11)}")
367         print(f"\t{'Cluster ID':>{len_id}} "
368                 f"| {'# of People':>{len_p_num}} "
369                 f"| {'Avg. of Ages':>{len_age}} "
370                 f"| {'Avg. of Severity':>{len_sev}} ")
371
372         # contents of table
373         cluster_id = 0
374         for people_num, avg_age, avg_sev in zip(people_num_of_a_cluster_list,
375                                                 avg_age_of_cluster_list,
376                                                 avg_severity_of_cluster_list):
377             print(f"\t{cluster_id:>{len_id}} "
378                     f"| {people_num:>{len_p_num}} "
379                     f"| {avg_age:>{len_age}} "
380                     f"|{avg_sev:>{len_sev}}")
381             cluster_id += 1
382
383         print(f"\t{'-'*(len_id+1)}"
384                 f"|{'-'*(len_p_num+2)}"
385                 f"|{'-'*(len_age+2)}"
386                 f"|{'-'*(len_sev+2)}-")
387
388         # bottom row
389         print(f"\t{'Total':^{len_id}} | {sum(people_num_of_a_cluster_list):>{len_p_num
                }} |")
```

```python
390         print(f"\t{'SSE':^{len_id}} | {round(self.sse_list[len(
            people_num_of_a_cluster_list) - 2], 2):>{len_p_num}} |")
391         print(f"\t{'Silhouette Score':>{len_id}} "
392             f"| {round(self.sil_score_list[len(people_num_of_a_cluster_list) - 2], 2
                ):>{len_p_num}} |")
393         print(f"\t{'-'*(len_sum+11)}")
394
395     def draw_silhouette(self):
396         """
397         method to draw graph using silhouette scores
398         :return: None
399         """
400         pass
401
402     def draw_graph(self):
403         """
404         method to draw clustering result
405         :return: None
406         """
407         pass
408
409     def draw_elbow_method(self, sse_list):
410         """
411         method to draw elbow graph using SSE(Sum of Squares Error)
412         :param sse_list: list of SSE
413         :return: None
414         """
415         plt.plot(range(2, 10), sse_list, marker='o')
416         plt.xlabel("The Number of Cluster")
417         plt.ylabel("SSE")
418         plt.show()
```